

# Prior distributions for variance-covariance matrices in hierarchical models

Edoardo Costantini

April 26, 2019

## 1 Introduction

Fully-Bayesian analysis has grown in popularity over the last few decades (look for reference). In particular, fully Bayesian approaches to hierarchical models are gaining momentum thanks to the wider accessibility of computing power (ehh.. verify!). The field is flourishing with both in terms of practical applications and theoretical interest.

Among the most active theoretical discussions, there is a long lasting debate on the value of subjective versus objective Bayesian perspectives (see Berger, 2006). While the former approach is often defended as the more truly "Bayesian", allowing the analyst to include prior information in the estimation process, the latter is appropriate when a researchers wants to express lack of subjective information on the model parameters of interest. The present work is not aiming to contribute to this debate directly, but rather is meant to focus on a particular issue pertaining to the objective framework.

When subscribing to the "objective" framework, a researcher can reflect the lack of prior information, relating to the parameters, by defining prior distributions that are meant to be minimally informative in some sense. Many theoretical contributions have tried to indicate what distributions, in which situations, are more apt to achieve such a goal. These endeavors have produced practical guidelines for researchers to assist them in the crucial decision prior distributions represent in Bayesian analysis.

In making this decision, two concepts are fundamental: conditional conjugacy and informativeness of a prior distribution. A family of priors for a parameter is conditionally conjugate when the conditional posterior of said parameter is also in that family of distributions. An uninformative prior is a reference prior distribution defined in such a way that the posterior inference is not influence in any way by it. Related to this concept, is that of a weakly informative prior. This is a proper distribution (it integrates to 1) but is set up so that "the information it does provide is intentionally weaker than whatever actual prior knowledge is available" (Gelman, 2006). This def-

initiation is rooted in the convenience, one may even say need, for a reference posterior distribution, by which is meant a distribution, obtained by employing an improper prior, that approximates a posterior that would have been obtained if a proper prior, describing initial vague information, had been used (Bernardo, 1979). In this paper, we will focus on proper priors, but the principle of uninformativeness still applies: objective Bayesian analysis needs a posterior distribution that does not incorporate the researcher’s personal beliefs so that it is possible to assess the relevance of initial information.

Much attention has been dedicated in the literature to propose prior distributions for scale parameters that are alternatives to the overused inverse Gamma distribution (reference dump here). This effort finds its *raison d’être* in the undesirable shrinkage of the posterior distributions towards zero, even when the parameters of an inverse Gamma prior are specified to achieve weak informativeness (Brown and Draper, 2006; Gelman, 2006). Many alternatives have been proposed but the focus has always been on hierarchical models with single random effects, even when the declared intent was to speak of cases of higher dimensionality (Kass and Natarajan, 2006)

My work contributes to this literature by exploring the degree of unformativity achieved by different prior distributions for the covariance matrix of a vector of random effects in a hierarchical model.

## 2 The hierarchical model of interest

I worked with a two-level normal model of repeated observations  $y_{ij}$  with individual (clusters) effects  $\mathbf{b}_i$

$$\begin{aligned} y_{ij} &= \mathbf{x}_{ij}^T \boldsymbol{\theta} + \mathbf{z}_{ij}^T \mathbf{b}_i + \epsilon_{ij} \\ \mathbf{b}_i &\sim N(\mathbf{0}, \boldsymbol{\Psi}) \\ \epsilon_{ij} &\sim N(0, \sigma^2) \end{aligned} \tag{1}$$

with  $\mathbf{b}_i$  a vector of random effects, and  $\boldsymbol{\Psi}$ , and  $\sigma^2$  representing the model hyperparameters. The model is fitted to a repeated observations data set where the main predictor is *time* ( $x_1$ ) and there is a dichotomous covariate  $x_2$  that interacts with it.  $\boldsymbol{\theta}$  is a vector of fixed effects containing fixed intercept, time, covariate, and interaction effects. This model is a random intercept and random slope model, with  $\mathbf{b}_i = [b_{i0}, b_{i1}]'$  a vector containing the cluster (individual) specific effects.

### 3 Weakly informative prior distributions in hierarchical models

The weak informativeness of a prior is inherently a provisional concept. It makes sense to fit a model with a specific prior definition, judge whether the posterior distribution makes sense (e.g. the values allowed by the posterior should be in a plausible range for the parameter of interest, and the posterior itself should not be shrieked to zero by the prior); and adjust the prior specification if the posterior does not make sense. In other words, if the posterior obtained using a given prior spreads over an unrealistic range of values for the parameter (e.g. excessively large variances), or if it is shrieked towards an implausible small number, then it is reasonable to reconsider the prior choice.

In the case of scalar random effects, such as for a random intercept model (as opposed to the vector of random effects that characterizes random intercept and slope models), it has been previously shown that the Inverse-gamma distribution with parameters  $\alpha = \epsilon$ ,  $\beta = \epsilon$ , commonly considered non informative, does not look to uphold to such a property (Gelman, 2006).

A fundamental issue of Bayesian Statistics is the definition of uninformative priors. It is common to distinguish between objective and subjective Bayesian approaches to the definition of priors (elaborate more). As Bayesian approaches to data analysis become more and more widespread a need for standard ways of defining priors for analyses becomes more and more relevant. While this issue might be easier for some types of analysis (elaborate more), when models become more complicated the issue becomes non-trivial. In particular there is no standard way of choosing an uninformative prior for the random effects variance-covariance matrix parameter in mixed effects models.

Publications have flourished on the matter but have often focused on the simpler case of a scalar random effect, the random intercept model. Here I will explain Brown and Draper 2006, Gelman 2006, Kass and Natarajan 2006 solutions. Some publications extended the reasoning to the vector random effects case, elaborating on the definition of a prior distribution for the variance covariance matrix (henceforth referred to as  $\Psi$ ), as random intercept and random slopes models.

Here I define the Bayesian model I want to use to test the performance of the different priors. The features I want to include are the following: continuous outcome, any number/measurement scale.

Let's start with the **model**

$$\begin{aligned} y_{ij} &= \mathbf{x}_{ij}^T \boldsymbol{\theta} + \mathbf{z}_{ij}^T \mathbf{b}_i + \epsilon_{ij} \\ \mathbf{b}_i &\sim N(\mathbf{0}, \Psi) \\ \epsilon_{ij} &\sim N(0, \sigma^2) \end{aligned} \tag{2}$$

(Vectors are in bold, matrix are capital Greek letters).

I'm going to define the following **priors**:

$$p(\boldsymbol{\theta}) \propto 1 \quad (3)$$

$$p(\sigma^2) \propto \sigma^{-2} \quad (4)$$

For what concerns the random effects variance covariance matrix, different priors are tested. In particular we used:

#### **inverse-Wishart prior**

$$p(\boldsymbol{\Psi}) \propto IW(\nu, S_0) \quad (5)$$

where  $\boldsymbol{\Psi}$  is the matrix of variance-covariance hyperparameters defined in model 1;  $S_0$  is a  $k \times k$  matrix, usually considered as a prior guess for the covariance matrix. The diagonal elements of  $\boldsymbol{\Psi}$  are inverse-Wishart distributed as well. In particular, the  $k_1 \times k_1$  upper-left triangular sub-matrix  $\boldsymbol{\Psi}_{11}$  is  $IW(\nu - (k - k_1), S_{011})$  distributed. Furthermore, note that in our case the diagonal elements are the random intercept and slope parameters, and therefore are scalar. For a univariate case, the inverse Wishart distribution simplifies to an inverse Gamma with parameters  $\alpha = \frac{\nu}{2}, \beta = \frac{S_{0kk}}{2}$ . This means that  $\sigma_{i0}^2 \sim IW(\nu - 1, S_{011})$  becomes  $\sigma_{i0}^2 \sim IG(\frac{\nu-1}{2}, \frac{S_{0ii}}{2})$ . By choosing  $\nu = k - 1 + \epsilon$  and  $S_{011} = k - 1 - 1 + \epsilon$  we obtain a prior on the variance parameters equivalent to the one  $IG(\epsilon, \epsilon)$  studied by Gelman (2006). Note that  $k$  being equal to two (there are two random effects) in model 1, the inverse Wishart prior for  $\boldsymbol{\Psi}$ , that grants an  $IG(\epsilon, \epsilon)$  on the variance parameters, is  $IW(1 + \epsilon, \epsilon I_2)$ .

- inverse-Wishart

$$p(\boldsymbol{\Psi}) \propto IW(\nu, S_0) \quad (6)$$

where we choose  $\nu = k - 1 + e$ , and  $S_0 = \text{diag}(k - 1 + e)$ , following indications by Gelman *et al* (2014). Given a  $\boldsymbol{\Psi} \sim IW(\nu, S_0)$ , where  $\boldsymbol{\Psi}$  and  $S_0$  are  $k \times k$  matrices, it is known that  $\boldsymbol{\Psi}_{11}$ , the  $k_1 \times k_1$  upper-left triangular sub-matrix of  $\boldsymbol{\Psi}$ , has an inverse-Wishart distribution as well. In particular,  $\boldsymbol{\Psi}_{11} \sim IW(\nu - (k - k_1), S_{011})$ . Furthermore, for a univariate case ( $k = 1$ ), we know that an inverse Wishart distribution simplifies to an inverse Gamma with parameters  $\alpha = \frac{\nu}{2}, \beta = \frac{S_{0kk}}{2}$ . With the goal of resembling as close as possible what Gelman 2006 did, we try to define inverse wishart priors for  $\boldsymbol{\Psi}$ , such that the the marginal disitbrution is as close as possible to the  $IG(e, e)$  used by Gelman. This goal is achieved by setting  $\nu = k - 1 + e$  and  $S_{011} = k - 1 + e$ , which makes the marginal distribution on the variance components (diagonals of  $\boldsymbol{\Psi}$ )  $\sigma_{ii}^2 \sim IG(\frac{\nu-1}{2}, \frac{S_{0ii}}{2})$ . The inverse-Wishart priors we defined are:

Prior Description	$\nu$	$S_0$
1. IW educated	2	educated guess
2. IW uninformative	$k - 1 + e$	$(k - 1 + e - 1) \times \text{diag}(k)$

1, .01, and .001 are then used as values of  $e$ .

- inverse-Wishart *a la* Huang and Wand

$$\begin{aligned}
p(\Psi|a_1, a_2) &\propto IW(\nu + k - 1, 2\nu \times \text{diag}(1/a_1, 1/a_2)), \\
a_k &\propto IG(1/2, 1/A_k^2),
\end{aligned} \tag{7}$$

with  $\nu = 2$  and  $\mathbf{A} = [1000, 1000]$ . The marginal distribution of any standard deviation term in  $\Psi$  is Half- $t(\nu, A_k)$  and, when choosing  $\nu = 2$ , the marginal distribution on the correlation term is uniform on  $(-1, 1)$ , see property 2 to 4 in Huang and Wand (2013, p. 442). Furthermore, according to Huang and Wand (2013, p. 441) arbitrarily large values for  $a_k$  lead to arbitrarily weak priors on the standard deviation term. Hence, our choices for the parameters of this prior are:

Prior Description	$\nu$	$\mathbf{A}$
3. IW a la HW	2	[1000, 1000]

- Matrix-F variate

Following Mulder, Pericchi 2018, I defined a matrix-F variate distribution as a prior for the covariance matrix of the random effects  $\Psi$

$$\begin{aligned}
p(\Psi) &\propto F(\Psi; \nu, \delta, \mathbf{B}) \\
&\propto \int IW(\Psi; \delta + k - 1, \Sigma) \times W(\Sigma; \nu, \mathbf{B}) d\Sigma
\end{aligned} \tag{8}$$

with degrees of freedom  $\nu > k - 1$ ,  $\delta > 0$ , and  $\mathbf{B}$  a positive definite scale matrix that functions as prior guess. Different strategies can be followed in trying to achieve vagueness of this prior. In the literature, the improper prior  $(\sigma^2)^{-\frac{1}{2}}$  has been proposed for the random effects variance (Berger, 2006; Berger and Strawderman, 1996). Placing a matrix-F prior on  $\Psi$ , one can approximate the improper  $|\Psi|^{-\frac{1}{2}}$  by choosing  $\nu = k$ , the smallest allowed integer,  $\mathbf{B} = b\mathbf{I}_k$  and letting  $b \rightarrow +\infty$  for fixed values of  $\delta$ .

Another approach might be trying to achieve a flat prior on the standard deviations of the random effects  $p(\sigma) \propto 1$ , through some other proper neighbor definition. Considering  $\Psi \sim \text{matrix-F}(\nu, \delta, \mathbf{B})$ , it is known that the marginal distribution of diagonal elements of  $\Psi$ ,

the variance components  $\sigma_{jj}^2$ , is univariate  $F(\nu, \delta, b_{jj})$  which is equivalent to  $p(\sigma_{jj}; \nu, \delta, b) \propto \sigma^{\nu-1} (1 + \sigma^2/b)^{-\frac{\nu+\delta}{2}}$ . By choosing  $\nu = k - 1 + \epsilon$ , with say  $\epsilon = .001$ ,  $\delta = 1$ , and  $b \rightarrow \infty$ ,  $p(\sigma) \propto 1$  is obtained. Hence, the specification  $\Psi \sim \text{matrix-F}(\nu = 1.001, \delta = 1, \mathbf{B} = 1e3\mathbf{I}_2)$  grants a proper prior neighbor to a flat prior on the standard deviation of the random effects.

It is also interesting to approach the vagueness issue by defining a vague prior weakly centered around an educated guess. Such educated guess can be obtained by plotting fitted regression lines to all the clusters in the data set that is under analysis, and trying to visually assess the intercepts and the slope variance. For example, consider the data set on depression that will be used to present the first empirical results. *here add plot and describe how you arrived to the educated guess* We can then try to achieve vagueness of the prior distribution by defining  $\nu$  as small as possible. Hence, another prior definition that was used in this paper, set  $\nu = k - 1 + \epsilon$ ,  $\delta = 1$  and  $\mathbf{B} = \text{matrix}(21, 0, 0, 9)$ .

Finally, we could also specify a prior that uses the empirical guess defined by Kass and Natarajan (2006). To do so, we specify the vague matrix-F prior with scale matrix  $\mathbf{R}^*$ , and  $\nu = 2$ .

The following table summarizes the matrix-F specifications that grant some form of non informativeness and that are compared in this paper.

Prior Description	$\nu$	$\delta$	$S_0$
$ \Psi ^{-\frac{1}{2}}$	2	1	$10^3 \times \mathbf{I}_2$
$p(\sigma) \propto 1$	1.001	1	$10^3 \times \mathbf{I}_2$
vaguely centered around educated guess	$k - 1 + \varepsilon$	$\varepsilon$	educated guess
vaguely centered around empirical $\mathbf{R}^*$	2	1	$\mathbf{R}^*$

$\varepsilon$  is set to each 1, and .1

## 4 Results

A good performance is shown by posteriors that give a reasonable range for the values of the parameters without having extremely long right-tails and without having excessively high peaks at zero.

In the figures, the empty circle represents the REML estimates of the parameter of interest. REML is used because it's the equivalent procedure to RIGLS, one of the most important ML estimation methods for Generalized Linear Models. RIGLS coincide with REML for all Gaussian models such as the one described in this article (Goldstein, 2010; Browne Draper, 2006).

When discussing the proper neighbor prior in the matrix F case for the SD of the variance, you should highlight that it is not an ideal situation, the role of a weakly informative prior is to regularize the posterior distribution so as to keep it into reasonable bounds (Gelman et al 2014

book). As a general guideline, priors should be made more precise as posteriors are more vague (which happens when fewer data are available).

When discussing the IW prior with e you can highlight that, compared to the uninformative matrix F priors, the range of values supported by the posterior distribution is decidedly more concentrated below 2 with a sharp peak towards 0.

In comparison, even when data is scarce (last few conditions), the posterior distributions obtained with the matrix F distribution allow for plausible values in a larger range that is more meaningful, while regularized by the prior to avoid impossible high values (which does happen when using the proper neighbor prior).

When discussing the priors for the correlation in the IW and matrix-F case, you should point out that the beta priors obtained, are weakly informative according to the symmetry principle: symmetrical prior distributions do not pull the posteriors in any particular direction.

## Derivation of Full Conditional Posterior Distribution

The derivation of the conditional posterior follows.

**Full conditional for  $\theta$**  (fixed effects)

Let's start with

$$p(\theta | \mathbf{y}, \mathbf{X}, \mathbf{Z}, \mathbf{b}_i, \Psi, \sigma^2) = p(\mathbf{y} | \theta, \mathbf{X}, \mathbf{Z}, \mathbf{b}_i, \Psi, \sigma^2) p(\theta)$$

where

$$\begin{aligned} p(\mathbf{y} | \theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) &= \prod_{i=1}^n \prod_{j=1}^J p(y_{ij} | \theta^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \Psi, \sigma^2) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} SSR\right) \end{aligned}$$

and

$$SSR = \sum_{i=1}^n \sum_{j=1}^J (y_{ij} - \theta^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2$$

where can rewrite  $y_{ij}$  as  $\tilde{y}_{ij}$ , with  $\tilde{y}_{ij} = y_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij}$  which makes SSR:

$$\begin{aligned} SSR &= \sum_{i=1}^n \sum_{j=1}^J (\tilde{y}_{ij} - \theta^T \mathbf{x}_{ij})^2 \\ &= (\tilde{\mathbf{y}} - \mathbf{X}\theta)^T (\tilde{\mathbf{y}} - \mathbf{X}\theta) \\ &= \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} - 2\theta^T \mathbf{X} \tilde{\mathbf{y}} + \theta^T \mathbf{X}^T \mathbf{X} \theta \end{aligned}$$

Hence,

$$p(\theta | \theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2} [-2\theta^T \mathbf{X} \tilde{\mathbf{y}} + \theta^T \mathbf{X}^T \mathbf{X} \theta]\right)$$

Combining this with the prior we obtain:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2) &\propto \exp\left(-\frac{1}{2}\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\theta}^T \mathbf{X} \tilde{\mathbf{y}}\right) \\ \boldsymbol{\theta}|\cdot &\sim N\left(\frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \tilde{\mathbf{y}}}{\sigma^2}, \frac{(\mathbf{X}^T \mathbf{X})^{-1}}{\sigma^2}\right) \end{aligned} \quad (9)$$

**Full conditional for  $\mathbf{b}_i$**  (random effects)

To derive this one we can start from:

$$p(\mathbf{b}_i|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\Psi}, \sigma^2) = p(\mathbf{y}_i|\boldsymbol{\theta}, \mathbf{b}_i, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2)p(\mathbf{b}_i)$$

We know that

$$p(\mathbf{y}_i|\cdot) = \prod_{j=1}^J p(y_{ij}|\boldsymbol{\theta}^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2} SSR_i\right)$$

with

$$SSR = \sum_{j=1}^J (y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2$$

and we can rewrite  $y_{ij}$  as  $\tilde{y}_{ij} = y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij}$ , which would make SSR be

$$\begin{aligned} SSR &= \sum_{j=1}^J (\tilde{y}_j - \boldsymbol{\theta}^T \mathbf{x}_j)^2 \\ &= (\tilde{\mathbf{y}} - \mathbf{b}_i^T \mathbf{Z}_i)^T (\tilde{\mathbf{y}} - \mathbf{b}_i^T \mathbf{Z}_i) \\ &= \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} - 2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i \end{aligned}$$

Hence,

$$p(\mathbf{y}_i|\cdot) \propto \exp\left(-\frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i]\right)$$

We also know that in this case, the "prior" is

$$p(\mathbf{b}_i) \propto N(\mathbf{0}, \boldsymbol{\Psi}) \propto \exp\left(-\frac{1}{2} [-2\mathbf{b}_i^T \boldsymbol{\Psi}^{-1} \mathbf{0} + \mathbf{b}_i^T \boldsymbol{\Psi}^{-1} \mathbf{b}_i]\right)$$

In conclusion, combining the sampling model and the prior, we get:

$$\begin{aligned} p(\mathbf{b}_i|\cdot) &\propto \exp\left(-\frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i] - \frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i]\right) \\ \mathbf{b}_i|\cdot &\propto N\left(\left(\boldsymbol{\Psi}^{-1} + \frac{\mathbf{Z}_i^T \mathbf{Z}_i}{\sigma^2}\right)^{-1} \left(\boldsymbol{\Psi}^{-1} \mathbf{0} + \frac{\mathbf{Z}_i^T \tilde{\mathbf{y}}_i}{\sigma^2}\right), \left(\boldsymbol{\Psi}^{-1} + \frac{\mathbf{Z}_i^T \mathbf{Z}_i}{\sigma^2}\right)^{-1}\right) \end{aligned} \quad (10)$$

**Full conditional for  $\sigma^2$**  (error variance)

The full conditional posterior can be expressed as:

$$p(\sigma^2|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{b}_i, \boldsymbol{\Psi}) = p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{b}_i, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2)p(\sigma^2)$$



The sampling model is the same we saw for the full conditional distribution of  $\theta$ :

$$\begin{aligned} p(\mathbf{y}|\theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) &= \prod_{i=1}^n \prod_{j=1}^J p(y_{ij}|\theta^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \Psi, \sigma^2) \\ &= \prod_{i=1}^n \prod_{j=1}^J (2\pi\sigma^{-2})^{-\frac{1}{2}} \exp\left(-\frac{(y_{ij} - \theta^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2}{2\sigma^2}\right) \end{aligned}$$

However, we are now interested in  $\sigma^2$ , hence

$$\begin{aligned} p(\mathbf{y}|\theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) &\propto (\sigma^2)^{-\frac{N}{2}} \exp\left(-\frac{\sum_{i=1}^n \sum_{j=1}^J (y_{ij} - \theta^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2}{2\sigma^2}\right) \\ &\propto (\sigma^2)^{-\frac{N}{2}} \exp\left(-\frac{1}{2\sigma^2} SSR\right) \end{aligned}$$

where  $N = \sum_i^n n_j$  is the entire sample size (all observations within all clusters). The prior for  $\sigma$  is given above, and therefore we can write the full conditional posterior as:

$$\begin{aligned} p(\sigma^2|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \theta, \mathbf{b}_i, \Psi) &\propto (\sigma^2)^{-\frac{N}{2}-1} \exp\left(-\frac{1}{2\sigma^2} SSR\right) \\ \sigma^2|. &\sim IG\left(\frac{N}{2}, \frac{SSR}{2}\right) \end{aligned} \tag{11}$$

**Full conditional for  $\Psi$**  (random effects variance covariance matrix)

Here, we need to write down the posteriors for the different priors we specified. First, let us define the sampling model for the random effects.

$$\begin{aligned} \begin{bmatrix} b_{0i} \\ b_{1i} \end{bmatrix} &= \mathbf{b}_i \sim N(\mathbf{0}, \Psi) \\ p(\mathbf{b}_1, \mathbf{b}_2|\Psi) &\propto |\Psi|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_b \Psi^{-1})\right) \end{aligned} \tag{12}$$

where  $\mathbf{S}_b$  is  $\sum_i \mathbf{b}_i \mathbf{b}_i^T$

- given the inverse-Wishart prior

$$\begin{aligned} p(\Psi) &\propto IW(\nu, \mathbf{S}_0) \\ &\propto |\Psi|^{-\frac{(\nu+n+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \Psi^{-1})\right) \end{aligned}$$

the full conditional posterior of  $\Psi$  is

$$\begin{aligned} p(\Psi|. ) &\propto |\Psi|^{-\frac{(\nu+n+k+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}([\mathbf{S}_0 + \mathbf{S}_b] \Psi^{-1})\right) \\ &\propto IW(\nu + n, \mathbf{S}_0 + \mathbf{S}_b) \end{aligned} \tag{13}$$

where  $\nu = 2$

- inverse-Wishart *a la* Huang and Wand

$$\begin{aligned}
p(\Psi|a_1, a_2) &\propto IW(\nu + k - 1, 2\nu \text{diag}(1/a_1, 1/a_2)), \\
a_k &\propto IG(\eta, 1/A_k^2) \\
p(\Psi) &\propto |\Psi|^{-\frac{(\nu+k-1+1)}{2}} \exp\left(-\frac{1}{2}\text{tr}(2\nu \text{diag}(1/a_1, 1/a_2)\Psi^{-1})\right) \\
&\quad \times \left(\frac{1}{a_1}\right)^{\eta+1} \exp\left(-\frac{1}{A_1^2 a_1}\right) \times \left(\frac{1}{a_2}\right)^{\eta+1} \exp\left(-\frac{1}{A_2^2 a_2}\right)
\end{aligned}$$

the full conditional posterior of  $\Psi$  is

$$\begin{aligned}
p(\Psi|\cdot) &\propto |\Psi|^{-\frac{(\nu+k-1+n+1)}{2}} \exp\left(-\frac{1}{2}\text{tr}([\mathbf{S}_b + 2\nu \text{diag}(1/a_1, 1/a_2)]\Psi^{-1})\right) \\
&\propto IW(\nu + k - 1 + n, \mathbf{S}_b + 2\nu \text{diag}(1/a_1, 1/a_2)) \\
p(a_k|\cdot) &\propto IG\left(\eta(\nu + k), \nu\left(\Psi_{kk}^{-1} + \frac{1}{A_k^2}\right)\right)
\end{aligned} \tag{14}$$

where  $\eta = \frac{1}{2}$ ,  $\nu = 2$ ,  $k = 2$ , and  $n$  is the number of clusters (individuals). (For the conditional posterior of  $a_k$  refer to Huang and Wand (2013), section 4.2).

- Matrix-F variate

Following section 2.3 in Mulder and Pericchi (2018), instead of working directly with the  $\Psi \sim F(\nu, \delta, \mathbf{B})$  we apply the parameter expansion defined above (see section on priors) and model it as  $\Psi \sim IW(\delta + k - 1, \mathbf{\Omega})$  with  $\mathbf{\Omega} \sim W(\nu, \mathbf{B})$ . With this parameter expansion, the conditional priors are:

$$\begin{aligned}
\Psi|\mathbf{\Omega} &\sim IW(\delta + k - 1, \mathbf{\Omega}) \\
\mathbf{\Omega}|\Psi &\sim W(\nu + \delta + k - 1, (\Psi^{-1} + \mathbf{B}^{-1})^{-1})
\end{aligned}$$

which makes the full conditional posterior of:

$$\begin{aligned}
\Psi|\mathbf{\Omega}, \cdot &\sim IW(\delta + k - 1 + n, \mathbf{S}_b + \mathbf{\Omega}) \\
\mathbf{\Omega}|\Psi, \cdot &\sim W(\nu + \delta + k - 1, (\Psi^{-1} + \mathbf{B}^{-1})^{-1})
\end{aligned}$$

with parameters as defined above. Given these posteriors, the Gibbs sampler implementation is straightforward.

## Notation Conventions

- $n$  number of clusters;  $i$  specific cluster
- $J$  number of observations within cluster;  $j$  specific observation
- $N$  total number of observations